

Трансформация подхода к хранению и синхронизации писем

Андрей Колесников



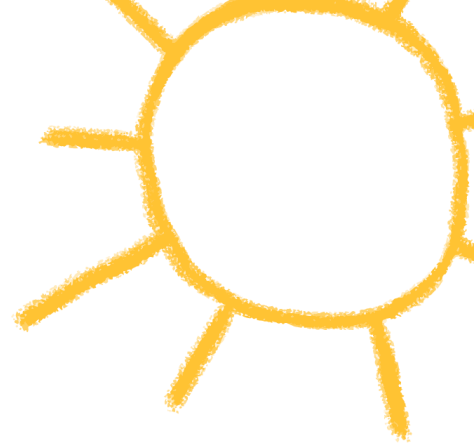
HighLoad⁺⁺
2022

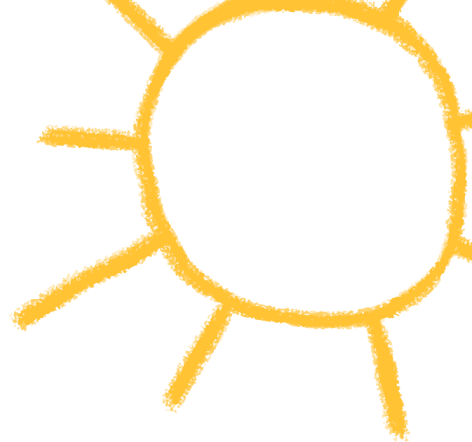
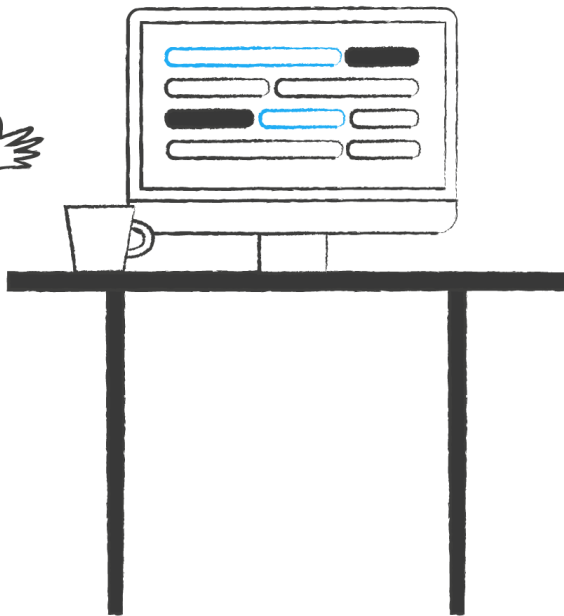
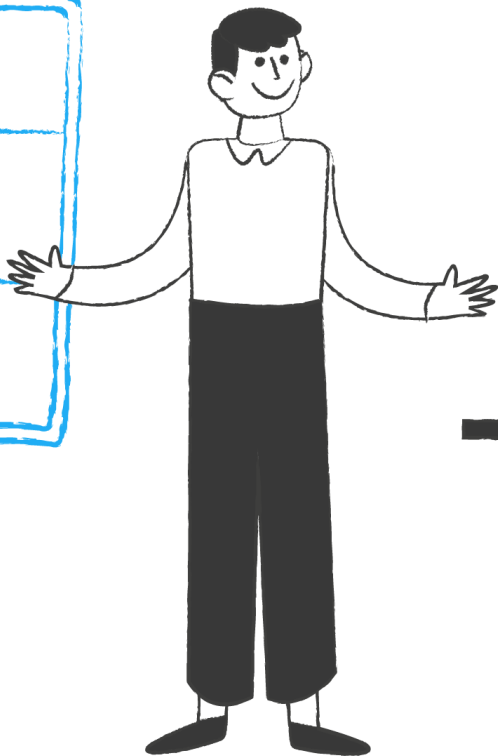
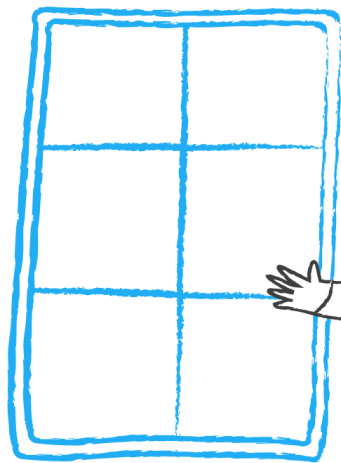
МойОфис

Колесников Андрей

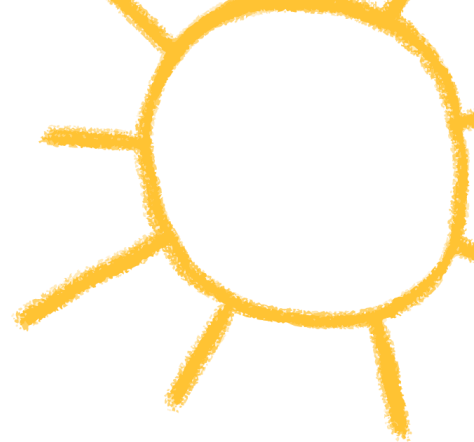
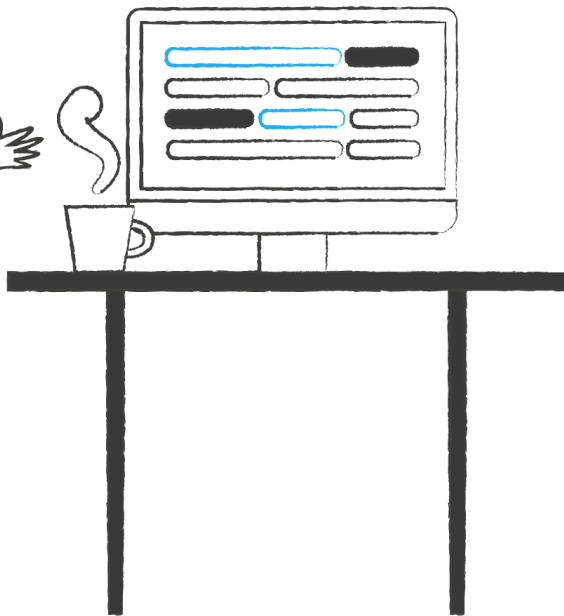
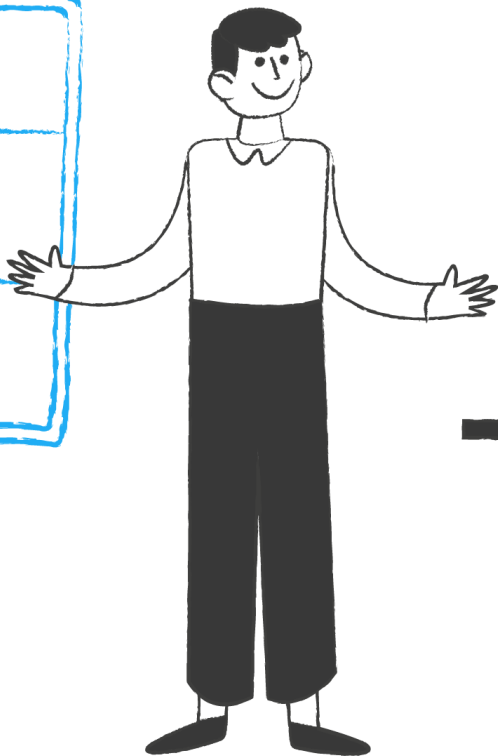
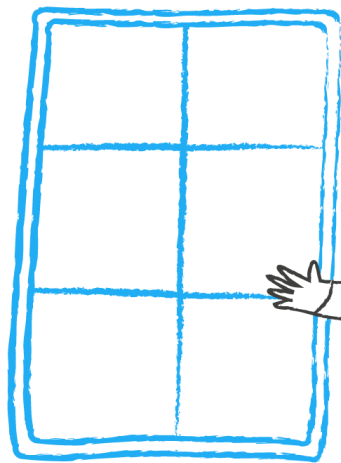
Руководитель инженерного отдела
МойОфис

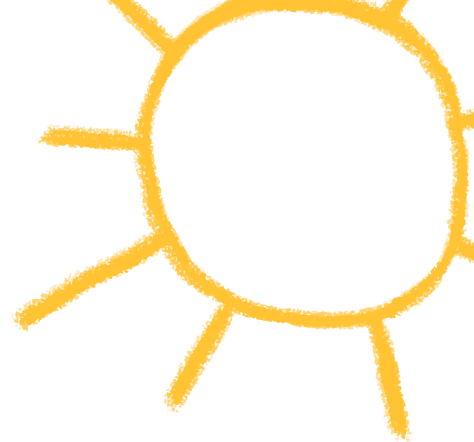
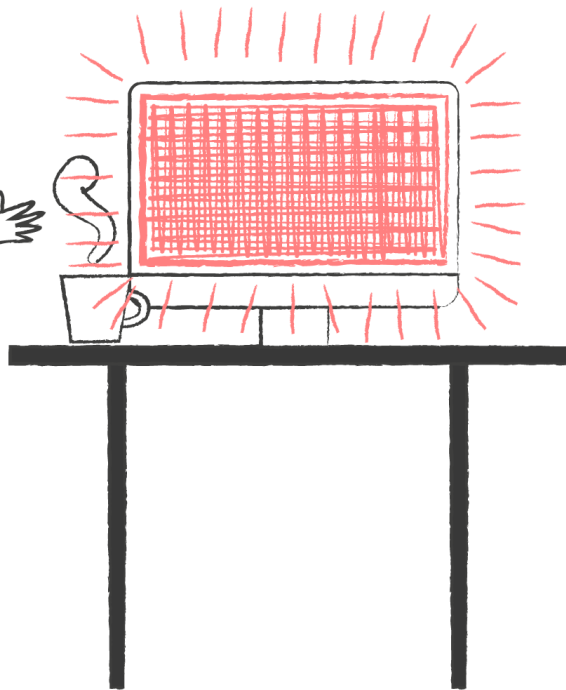
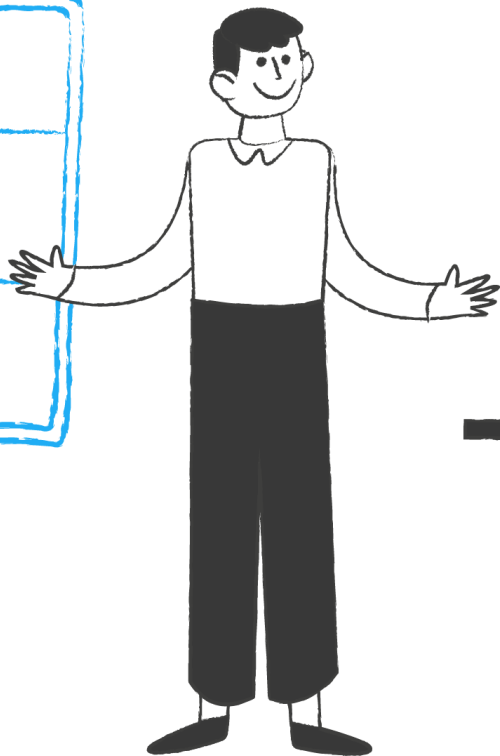
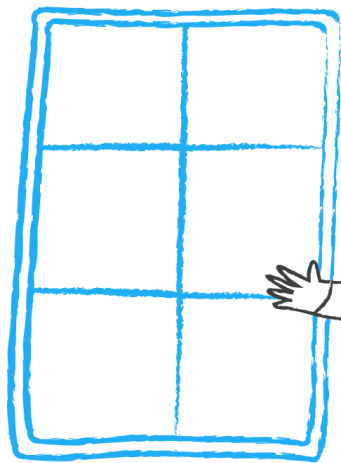


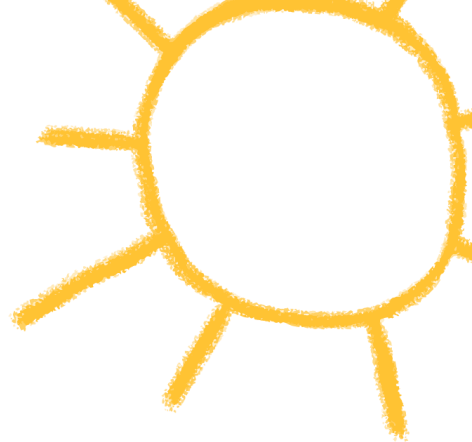
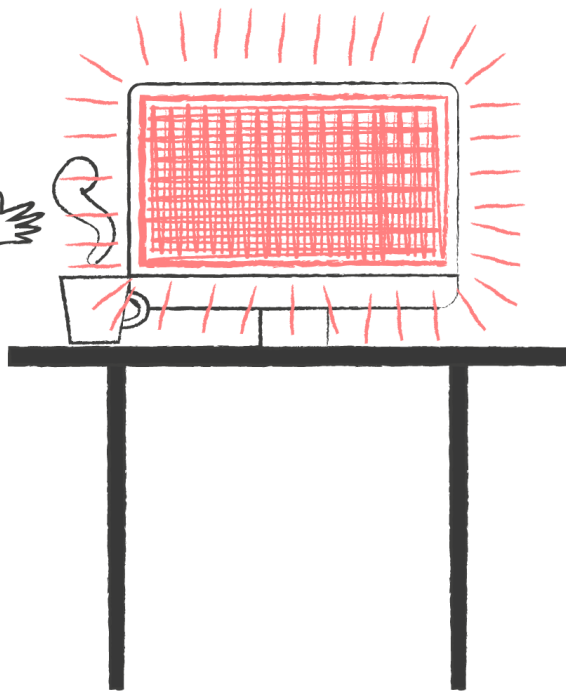
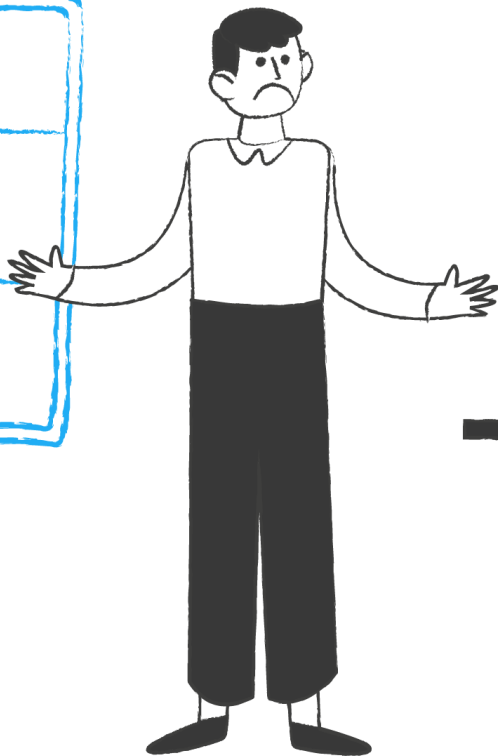
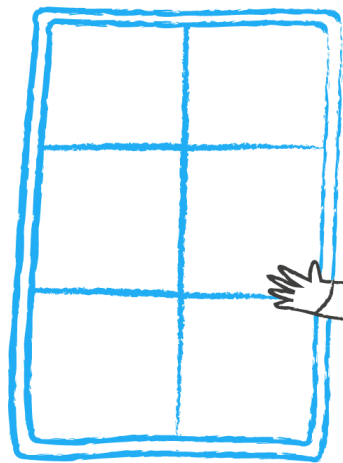




МойОфис

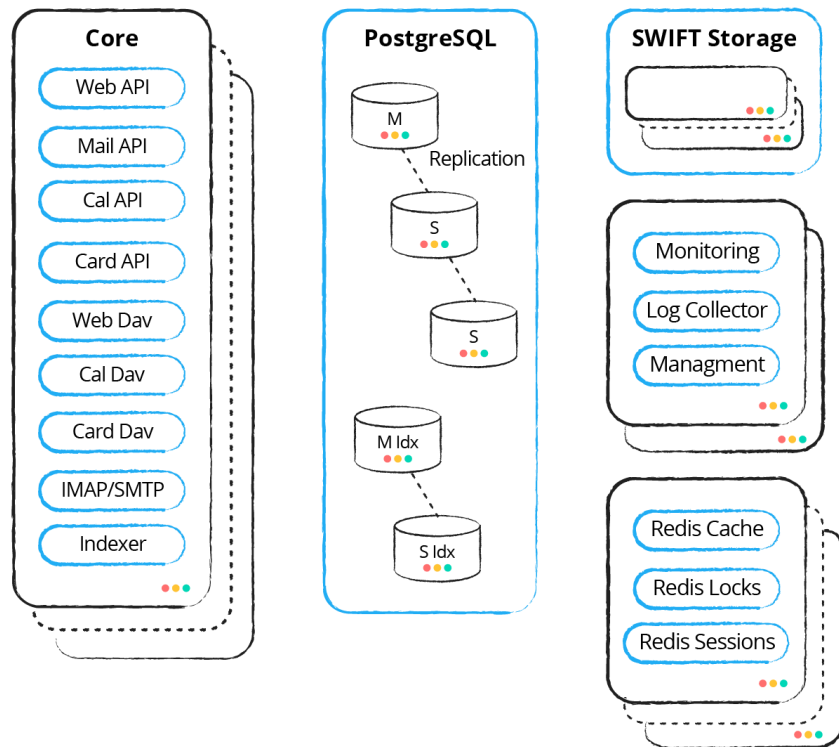






Начало: PostgreSQL, OpenStack Swift

- Метаданные хранятся в PostgreSQL (конверт), тело — в OpenStack Swift
- Кластеризация Расетmaker (менеджер ресурсов) + Corosync (транспортный уровень)



Начало: PostgreSQL, OpenStack Swift

Посмотреть заголовок в базе:

```
SELECT * FROM objects WHERE props->>'envelope' LIKE '%F36A916005E@myoffice.ru%';
```

```
-----  
id      | Oalfa1hDynrh3tNhluiB  
owner   | Ualfa1TT8sOwtqAoMo6  
...  
title   | Совещание:1 – тема сообщения  
...  
p_file  | Falfa1hDynrh6axzBS8J
```

Начало: PostgreSQL, OpenStack Swift

Посмотреть заголовок в базе:

```
SELECT * FROM objects WHERE props->>'envelope' LIKE '%F36A916005E@myoffice.ru%';
```

```
-----  
id      | Oalfa1hDynrh3tNhluiB  
owner   | Ualfa1TT8sOwtqAoMo6  
...  
title   | Совещание:1 – тема сообщения  
...  
p_file  | Falfa1hDynrh6axzBS8J
```

Получить файл письма:

```
swift -A http://10.2.3.4:8280/auth/v1.0 -U system:syncacc -K XXXXXXXXXXXX download  
Ualfa1TT8sOwtqAoMo6 Falfa1hDynrh6axzBS8J
```

Начало: PostgreSQL, OpenStack Swift

Посмотреть заголовок в базе:

```
SELECT * FROM objects WHERE props->>'envelope' LIKE '%F36A916005E@myoffice.ru%';
```

```
-----  
id          | Oalfa1hDynrh3tNhluiB  
owner       | Ualfa1TT8sOwtqAoMo6  
...  
title       | Совещание:1 – тема сообщения  
...  
p_file      | Falfa1hDynrh6axzBS8J
```

Получить файл письма:

```
swift -A http://10.2.3.4:8280/auth/v1.0 -U system:syncacc -K XXXXXXXXXXXX download  
Ualfa1TT8sOwtqAoMo6 Falfa1hDynrh6axzBS8J
```

```
more Falfa1hDynrh6axzBS8J
```

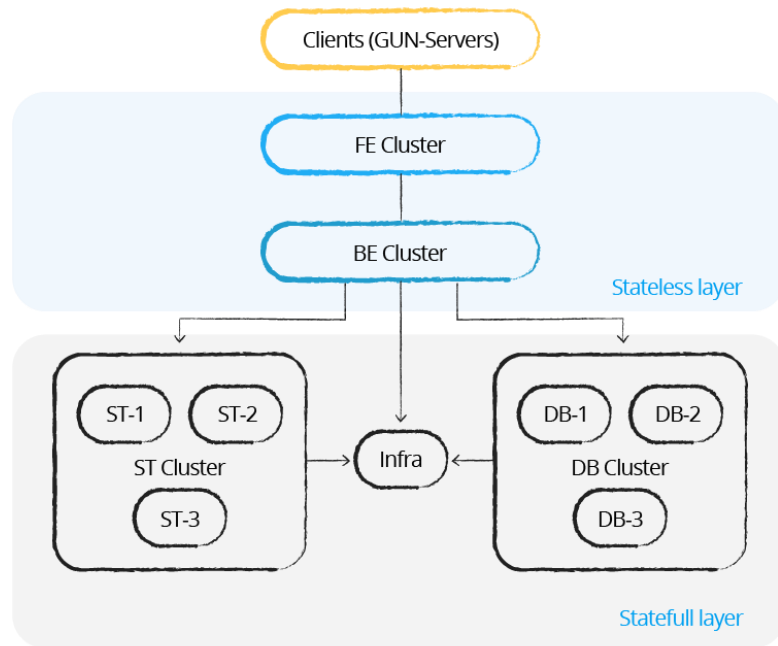


Предел производительности

Максимальная производительность
2500 RPS (100 000 пользователей)

За время тестирования в системе
было создано:

- 342 073 пользователя
- 13 833 934 директории
- 18 628 870 файлов



PostgreSQL, OpenStack Swift



Плюсы:

- Наличие экспертизы
- Проверенные временем технологии
- Большое комьюнити

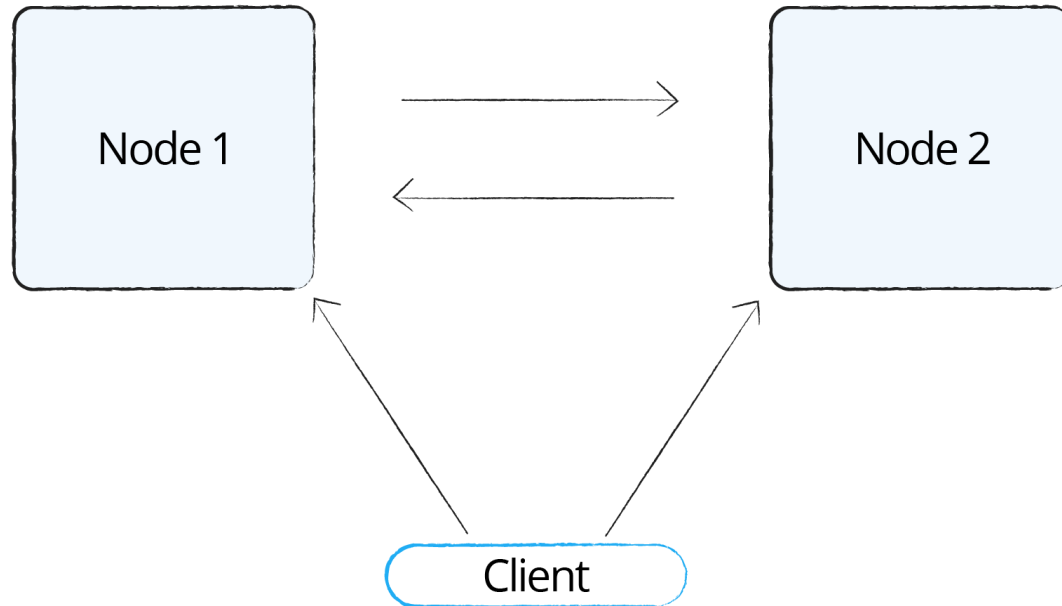


Минусы:

- Архитектурные погрешности, монолит
- Ограничения масштабирования (PostgreSQL по умолчанию поставляется для работы с одним сервером, обслуживающим запросы и двумя резервными)
- Подсистема управления блокировками и кэшированием (Redis) — узкое место для большинства методов

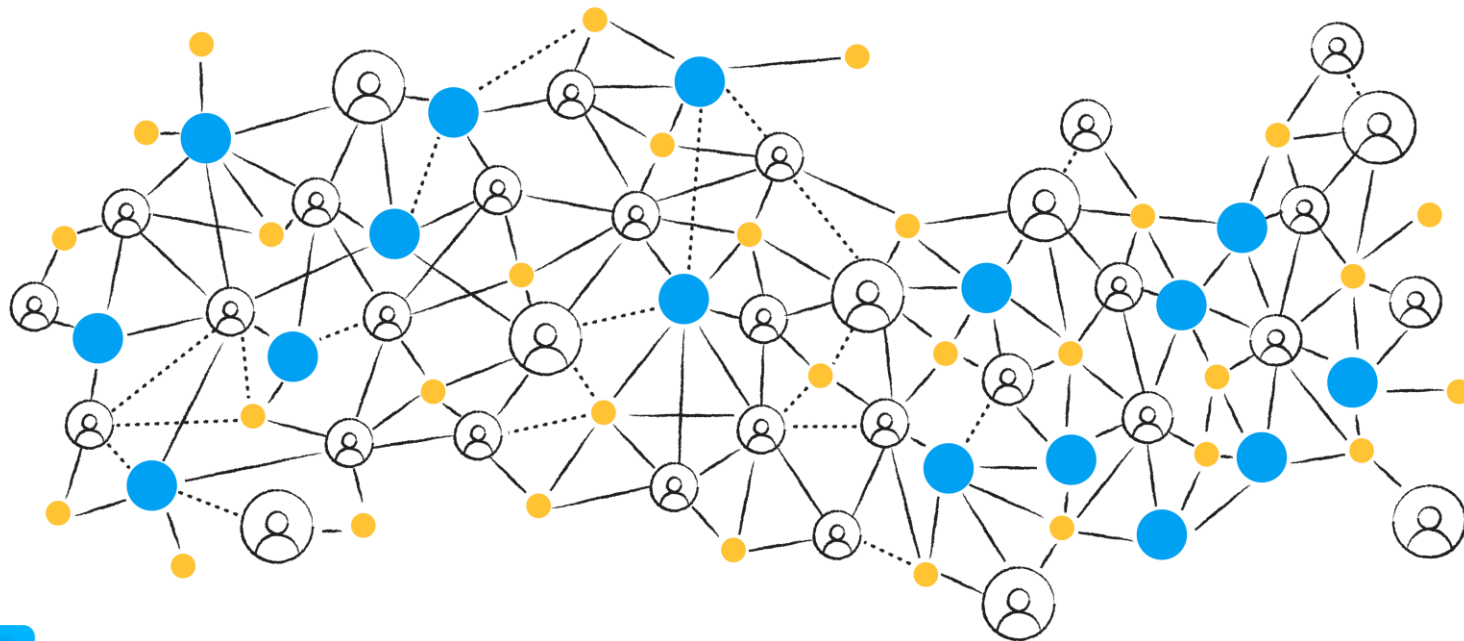


Dovecot dsync



Dovecot dsync & Shared Mailboxes

- Расшаренные ящики == проблемы синхронизации dsync:
<https://wiki.dovecot.org/SharedMailboxes/ClusterSetup>

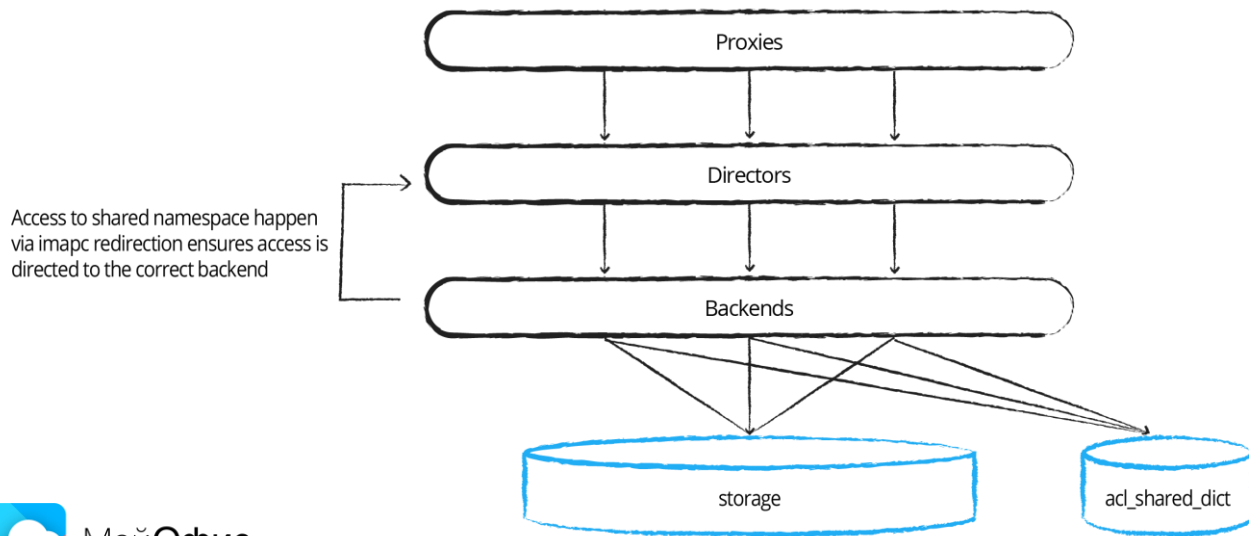


Dovecot dsync: incident resolved

<https://wiki.dovecot.org/SharedMailboxes/ClusterSetup>

There are some limitations for this kind of use case:

- `imapc_*` settings are global. You can't have two different namespaces with different `imapc` setting yet.
- The `imapc` code doesn't support some IMAP features. Most importantly `SORT` isn't supported, which may result in lower performance.



GlusterFS

Brick – каталог хранения, основная единица хранения GlusterFS

Volume – том, логическая комбинация каталогов хранилища

```
gluster volume info
```

```
Volume Name: gv0
```

```
Type: Replicate
```

```
Volume ID: f25cc3d8-631f-41bd-96e1-3e22a4c6f71f
```

```
Status: Started
```

```
Snapshot Count: 0
```

```
Number of Bricks: 1 x 3 = 3
```

```
Transport-type: tcp
```

```
Bricks:
```

```
Brick1: server1:/data/brick1/gv0
```

```
Brick2: server2:/data/brick1/gv0
```

```
Brick3: server3:/data/brick1/gv0
```

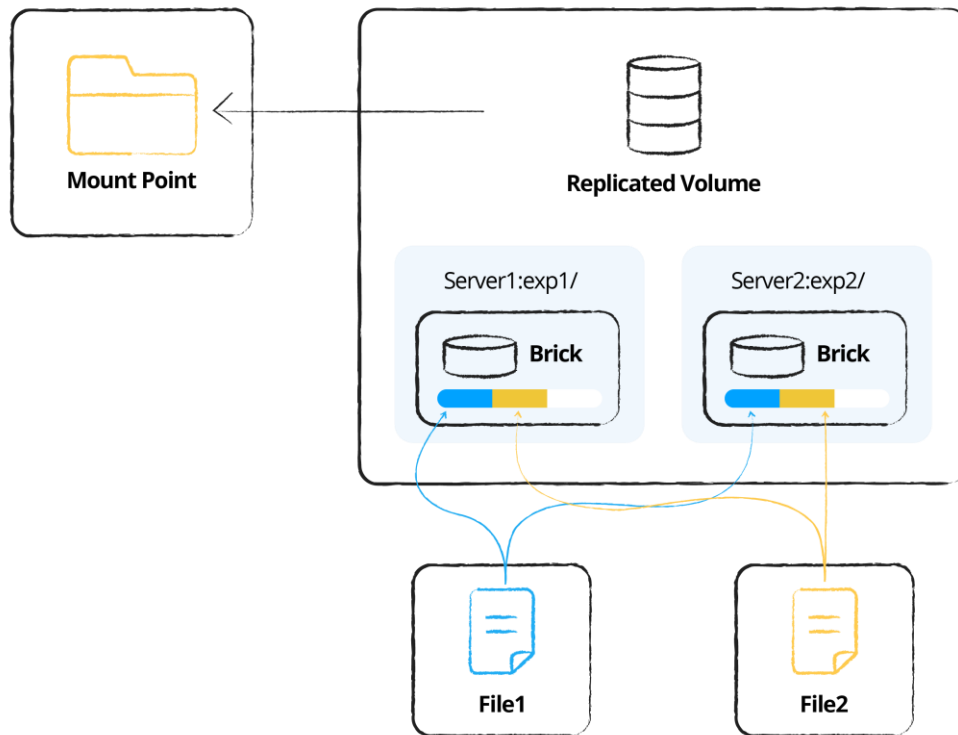
```
Options Reconfigured:
```

```
transport.address-family: inet
```

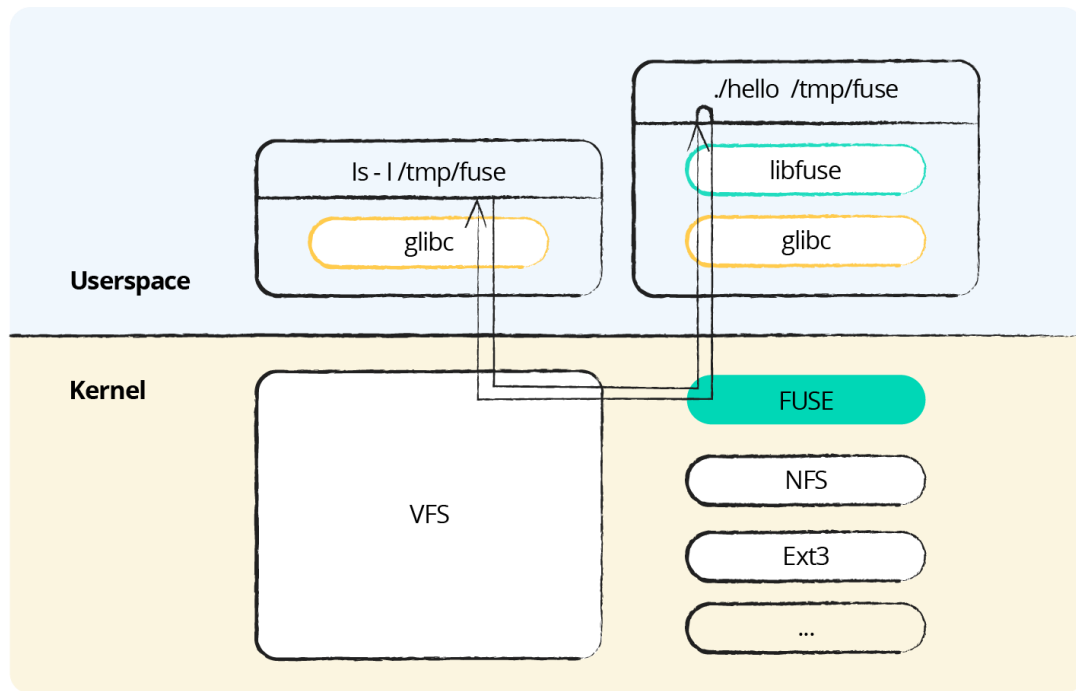


GlusterFS

Replicated — тома с репликацией. Аналогично RAID 1. В такой конфигурации одни и те же данные записываются минимум на два подтома.



GlusterFS



GlusterFS



Плюсы:

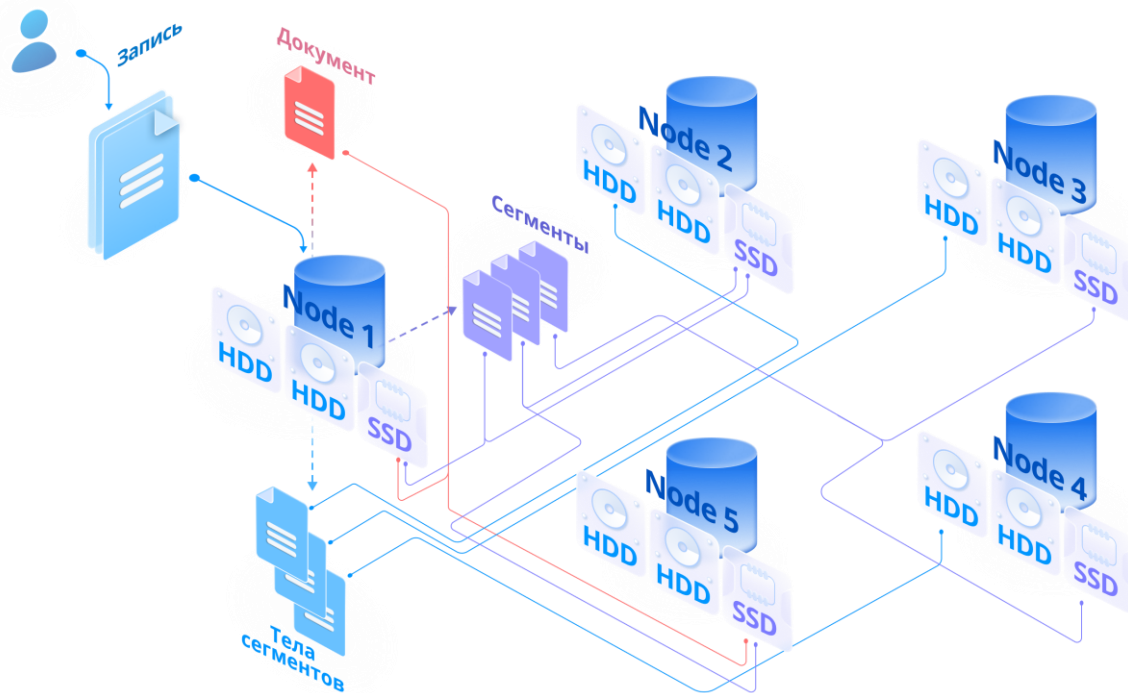
- Быстрый старт
- Большое комьюнити и сопровождение Red Hat (В 2011 проект был приобретен компанией Red Hat и лег в основу Red Hat Storage Server)



Минусы:

- Чувствительность к задержкам в сети ($> 50\text{ ms}$)
- Массовое удаление файлов: Remote I/O error

Собственная разработка: Dispersed Object Store



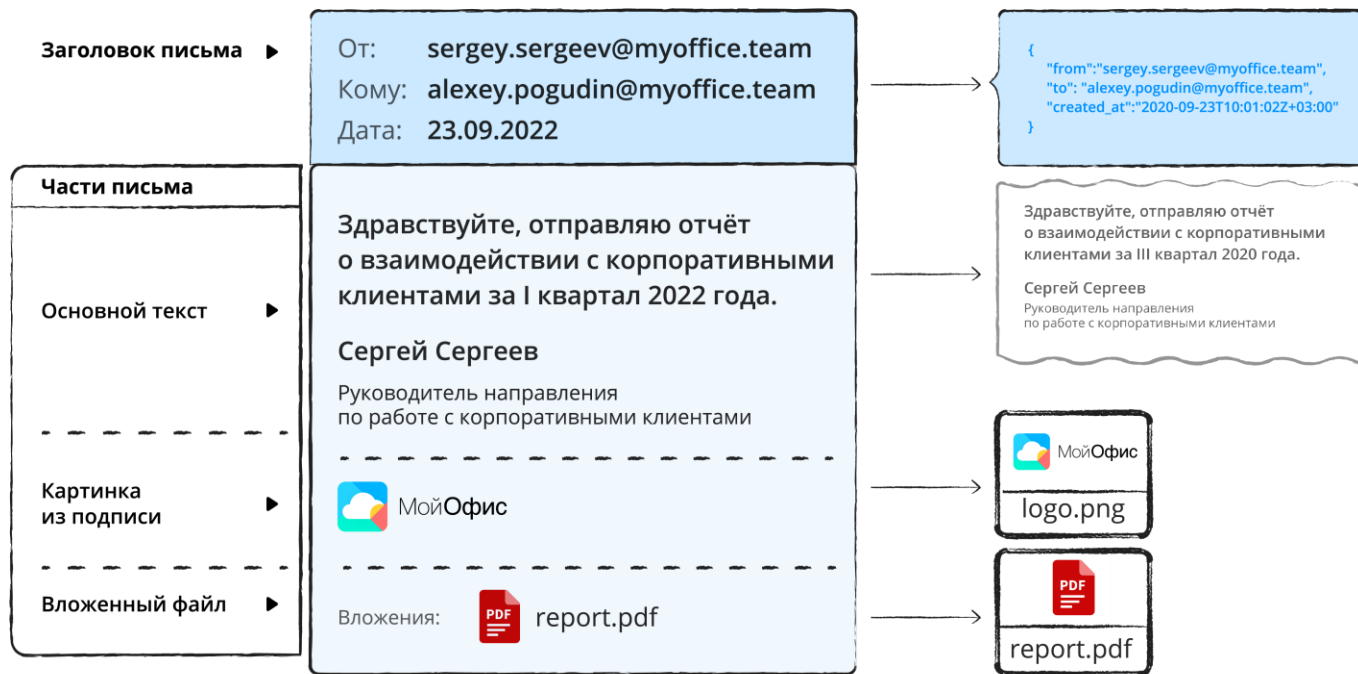
Особенности корпоративной переписки



**Большое количество полностью
или частично совпадающих писем:**

- У каждого письма минимум 2 копии (отправитель и получатель).
- Корпоративные рассылки (ещё больше копий).
- Шаблонные письма (ERP-системы, баг-трекеры, календари).
- Повторяющиеся элементы писем (вложения, подписи).
- Длинные цепочки писем со взаимным цитированием.

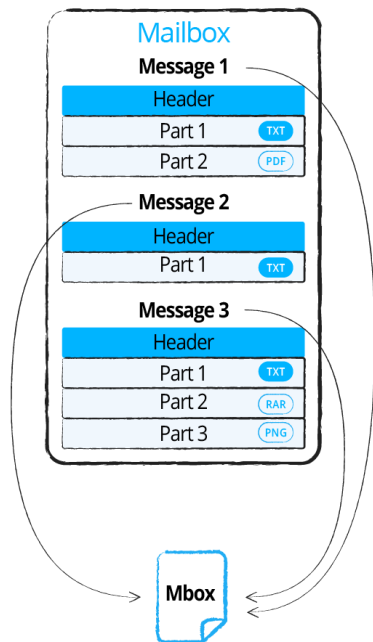
Структура электронного письма



Способы хранения электронной почты

Mbox

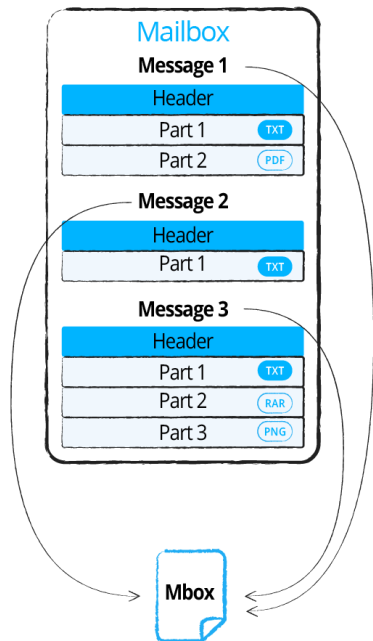
один почтовый ящик — один файл



Способы хранения электронной почты

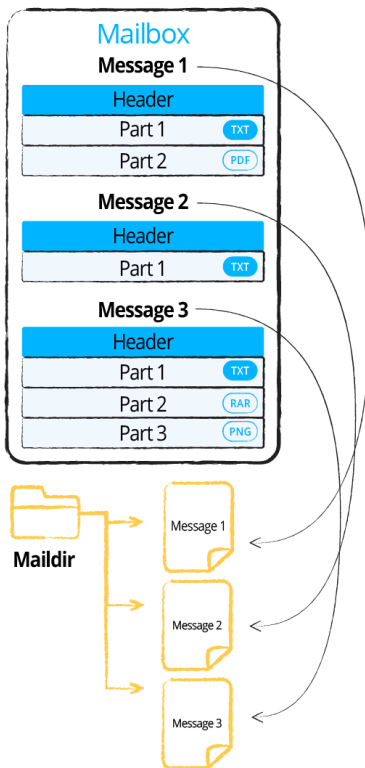
Mbox

один почтовый ящик — один файл



Maildir

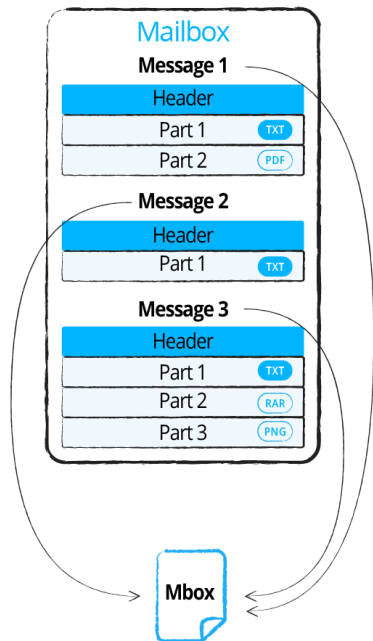
одно письмо — один файл



Способы хранения электронной почты

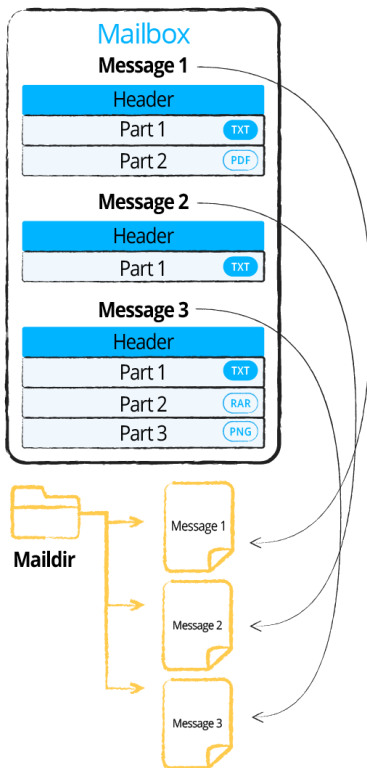
Mbox

один почтовый ящик — один файл



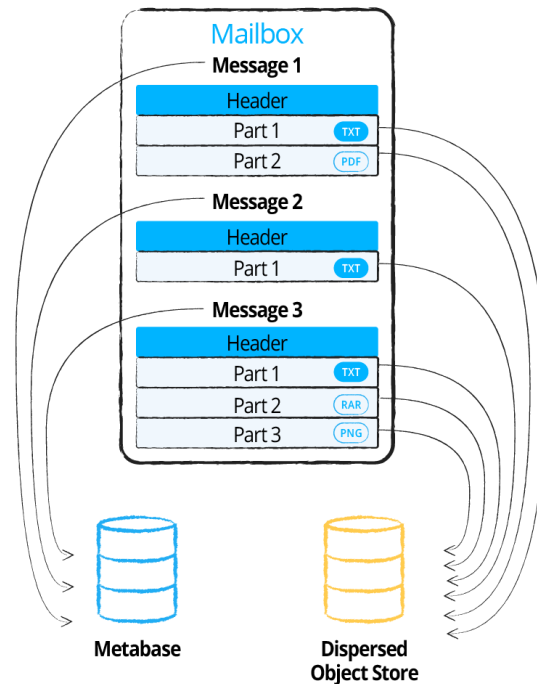
Maildir

одно письмо — один файл



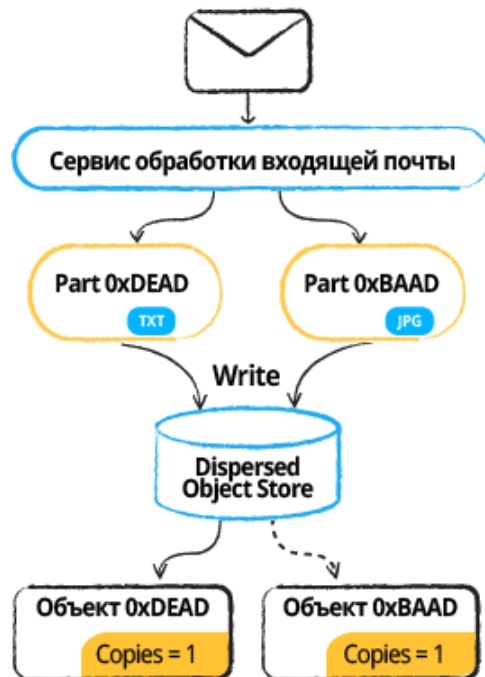
Mailon

один парт — один объект



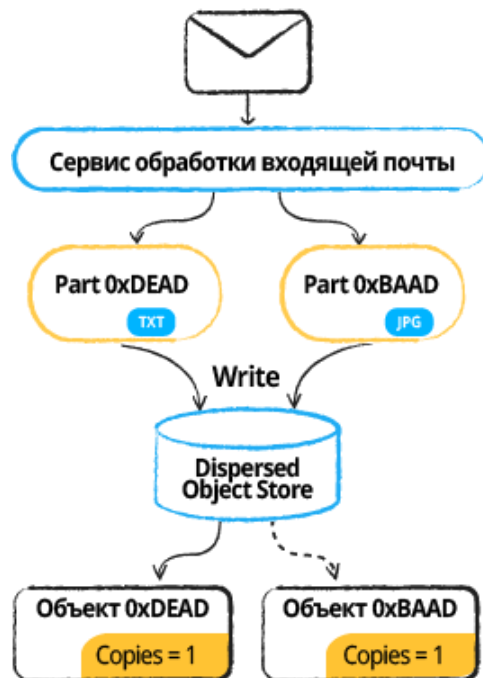
Дедупликация на уровне объектов

Шаг №1: сохранение
первого письма

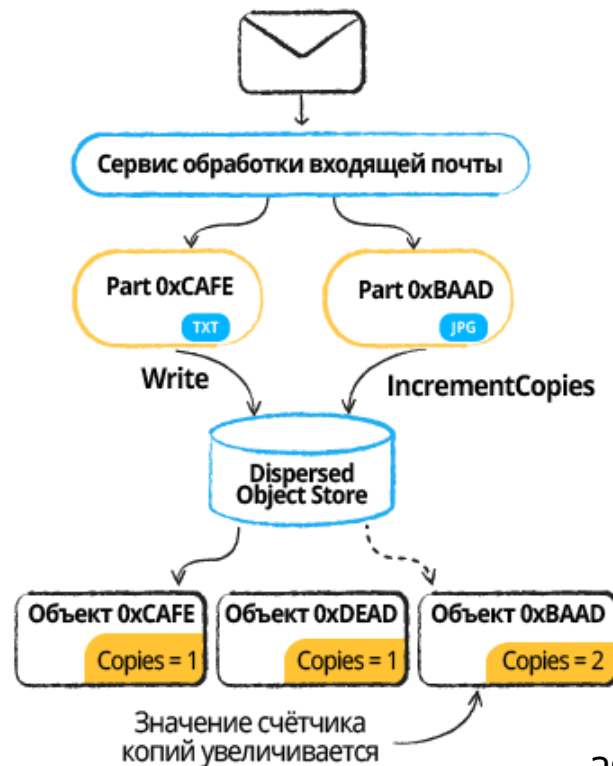


Дедупликация на уровне объектов

Шаг №1: сохранение первого письма



Шаг №2: сохранение второго письма, которое частично совпадает с первым



Дедупликация на уровне объектов



Плюсы:

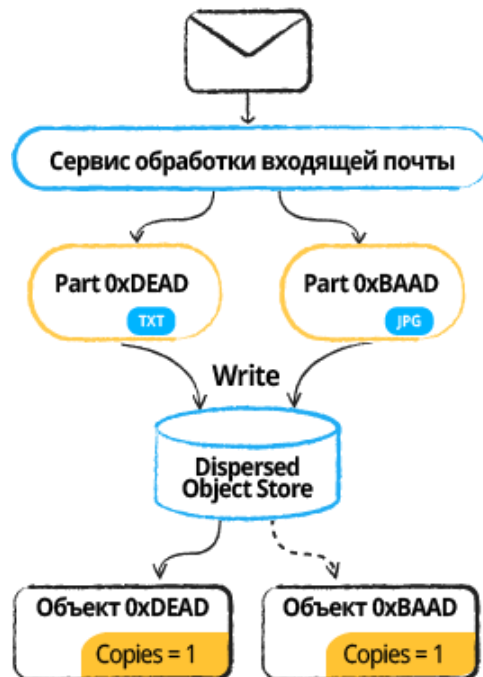
- Экономия дискового пространства
- Экономия сетевого трафика



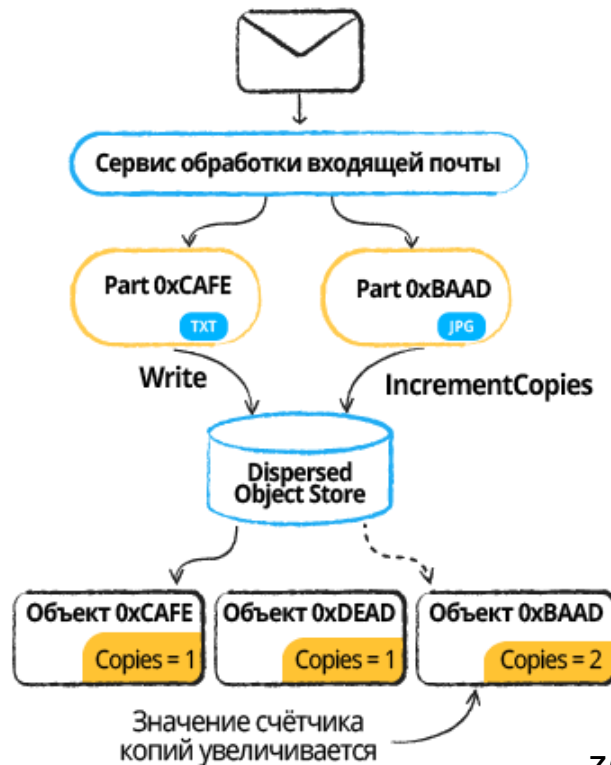
Минусы:

- Усложнение логики клиентского сервиса

Шаг №1: сохранение первого письма



Шаг №2: сохранение второго письма, которое частично совпадает с первым



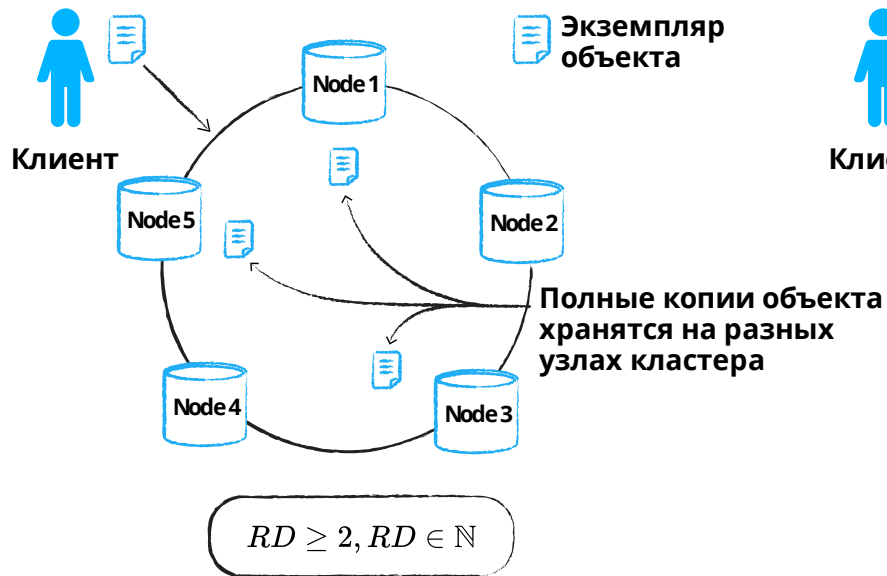
Тип и размер данных

- Хранилище умеет определять формат поступивших данных:
 - ✓ Text (txt, html, rfc822)
 - ✓ Binary (image, video, audio, pdf)
- Разные конвейеры обработки для разных типов данных.
- Размер данных имеет значение.

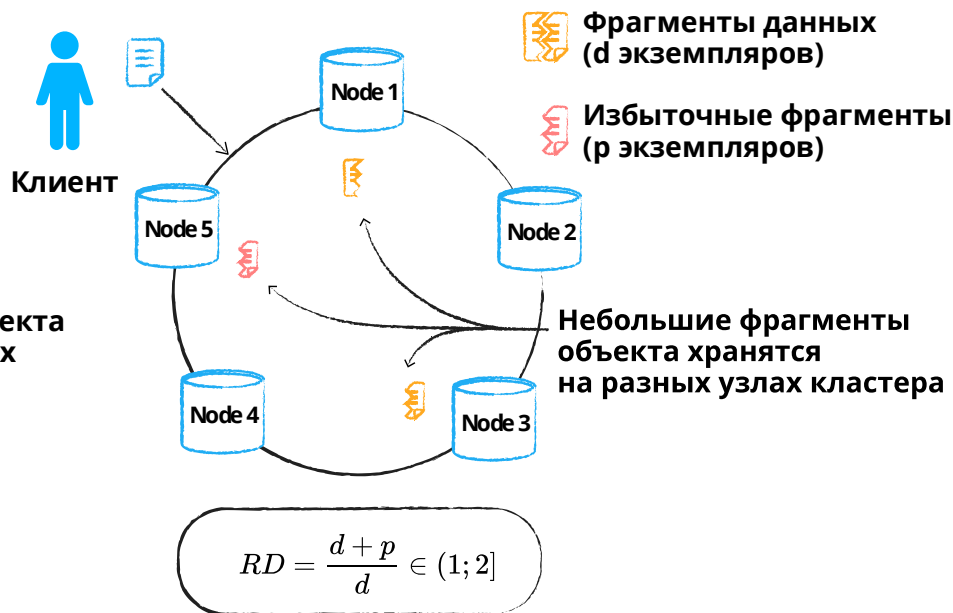
Класс размеров	Текстовые данные		Бинарные данные	
	Чанкинг	Компрессия	Чанкинг	Компрессия
Большие, средние	Content-defined chunking	+	Equal-size chunking	-
Малые	-	-	-	-

Отказоустойчивость и избыточность

«Полная» избыточность



«Дробная» избыточность



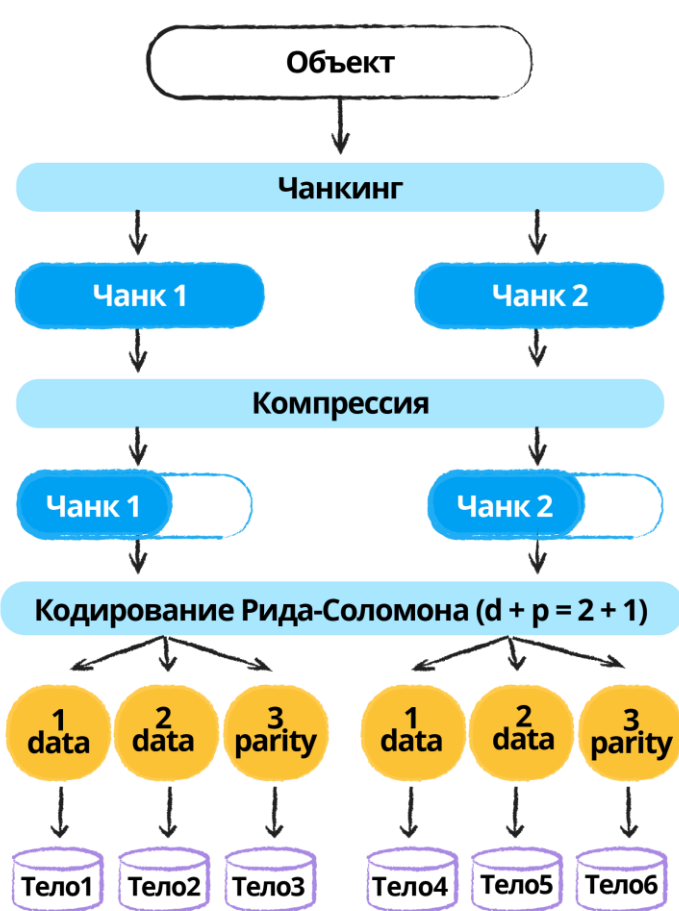
Отказоустойчивость и избыточность

- Соотношение метаданных к данным изменяется в пределах 1:3 — 1:500.
- Для каждого слоя хранения установлен свой уровень избыточности.

$$RD_{meta} \in \{2, 3, 5\}$$
$$RD_{data} \in (1; 2]$$

Слой хранения метаданных (RF копии)

Слой хранения данных (1 копия)



Эффективность оптимизаций



RECEIVED

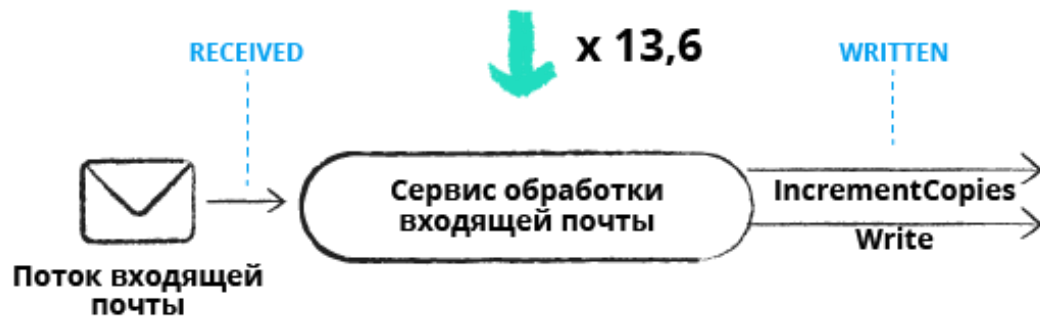


$$RD_{meta} = 3$$
$$RD_{data} = \frac{2+1}{2} = 1.5$$



МойОфис

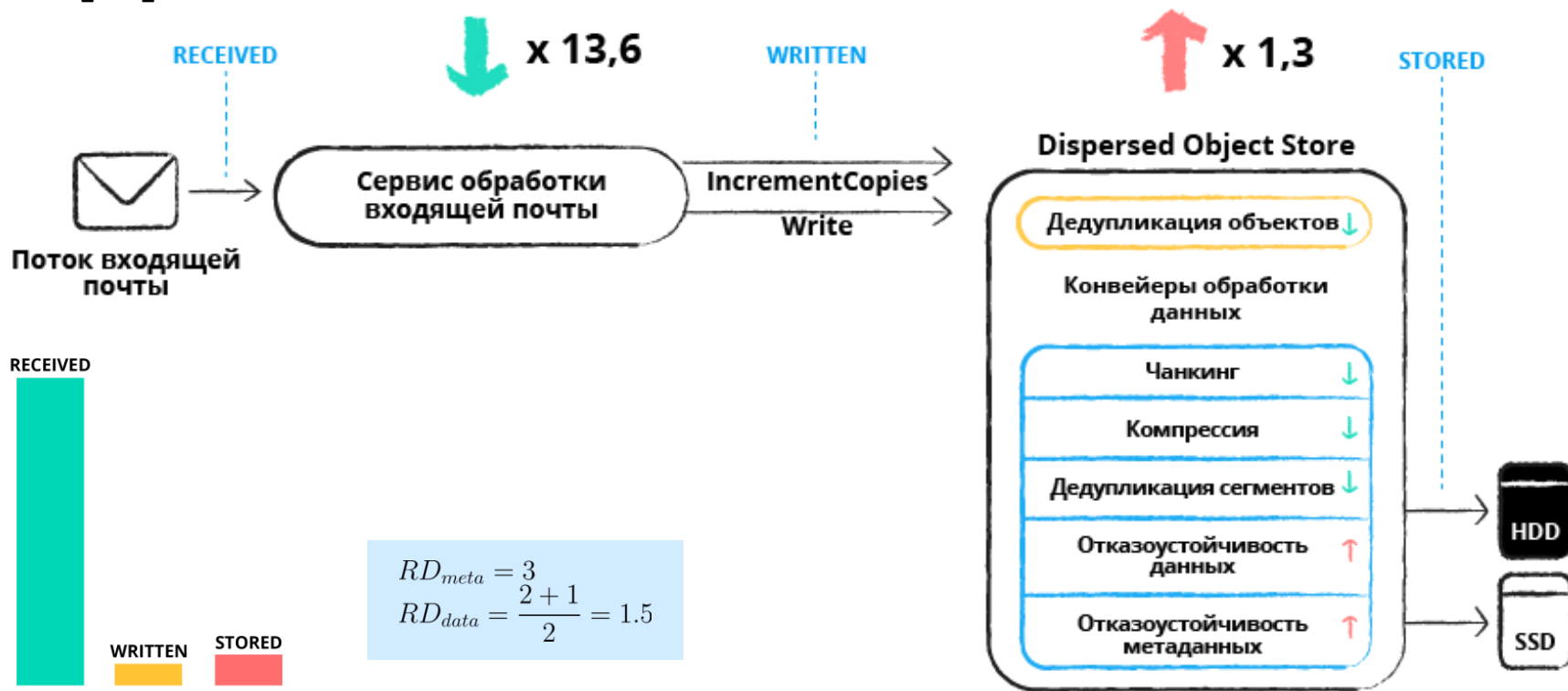
Эффективность оптимизаций



$$RD_{meta} = 3$$
$$RD_{data} = \frac{2+1}{2} = 1.5$$



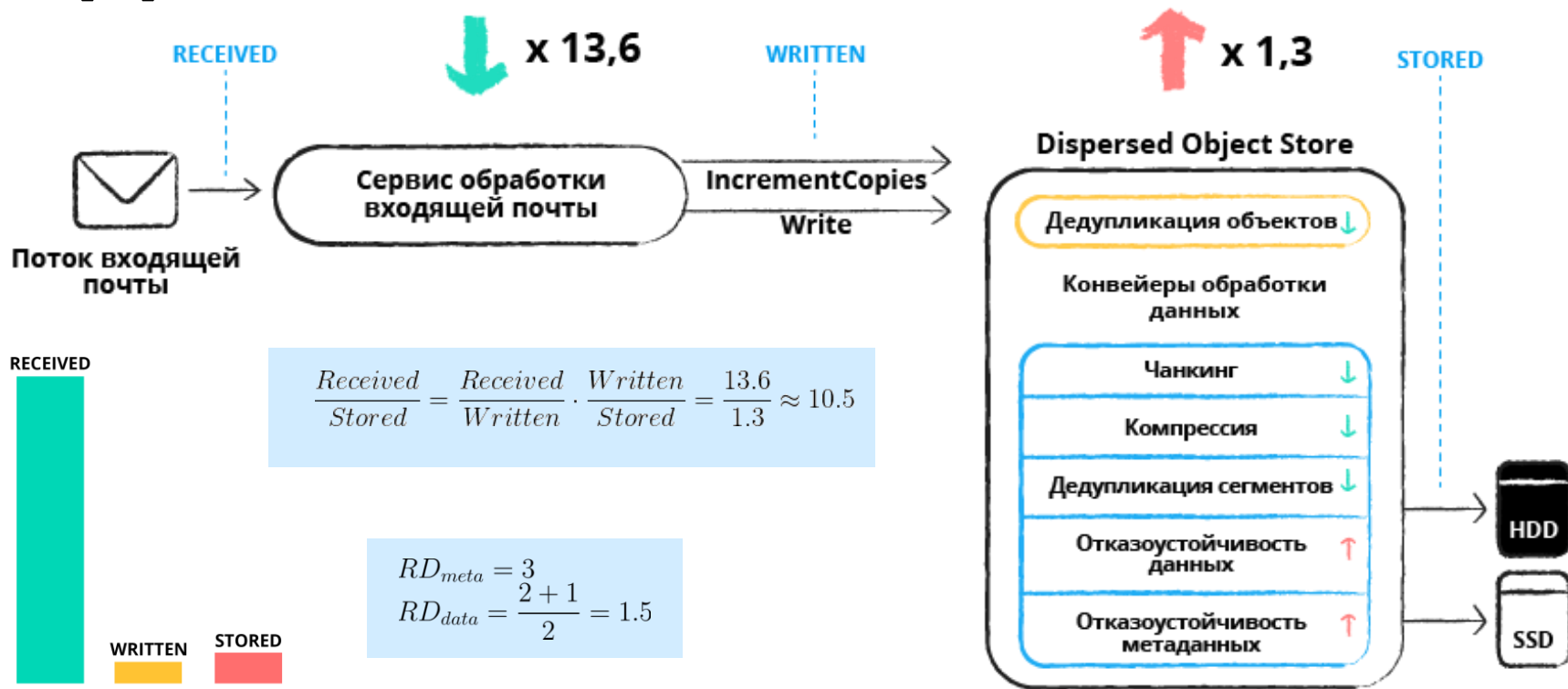
Эффективность оптимизаций



$$RD_{meta} = 3$$
$$RD_{data} = \frac{2+1}{2} = 1.5$$

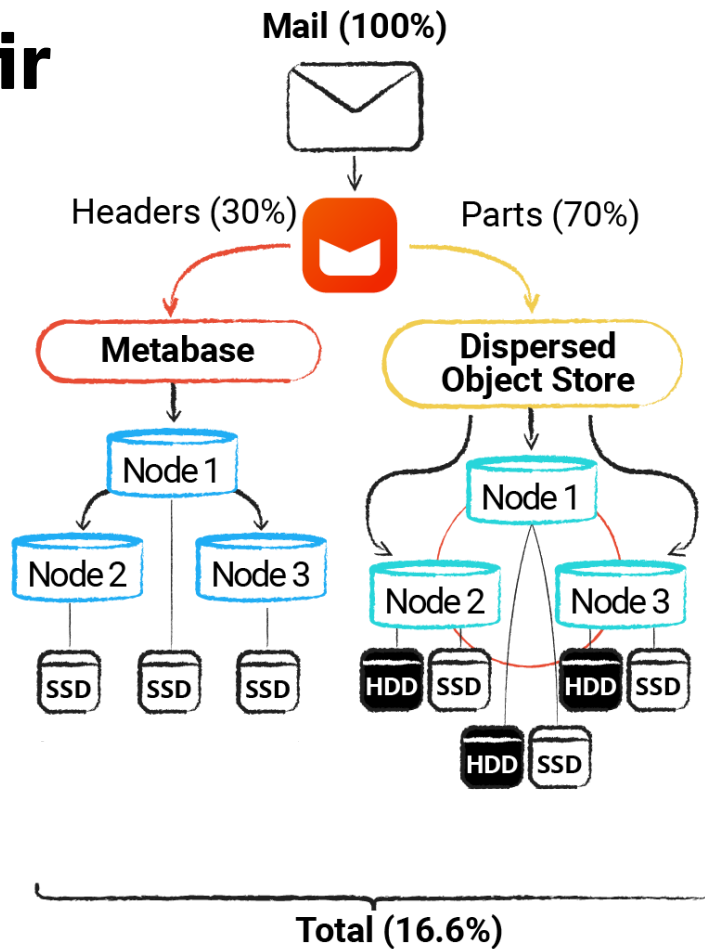


Эффективность оптимизаций



Mailion vs Maildir

$$RD_{meta} = 3$$
$$RD_{data} = \frac{2+2}{2} = 2$$

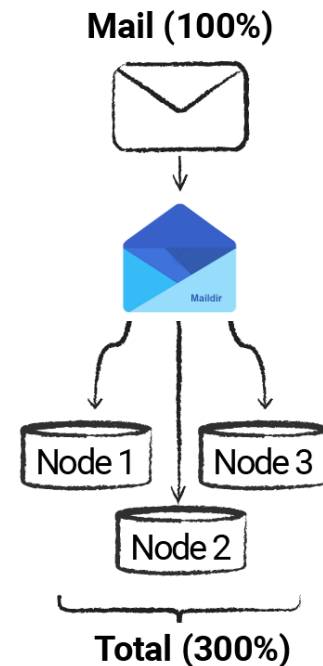
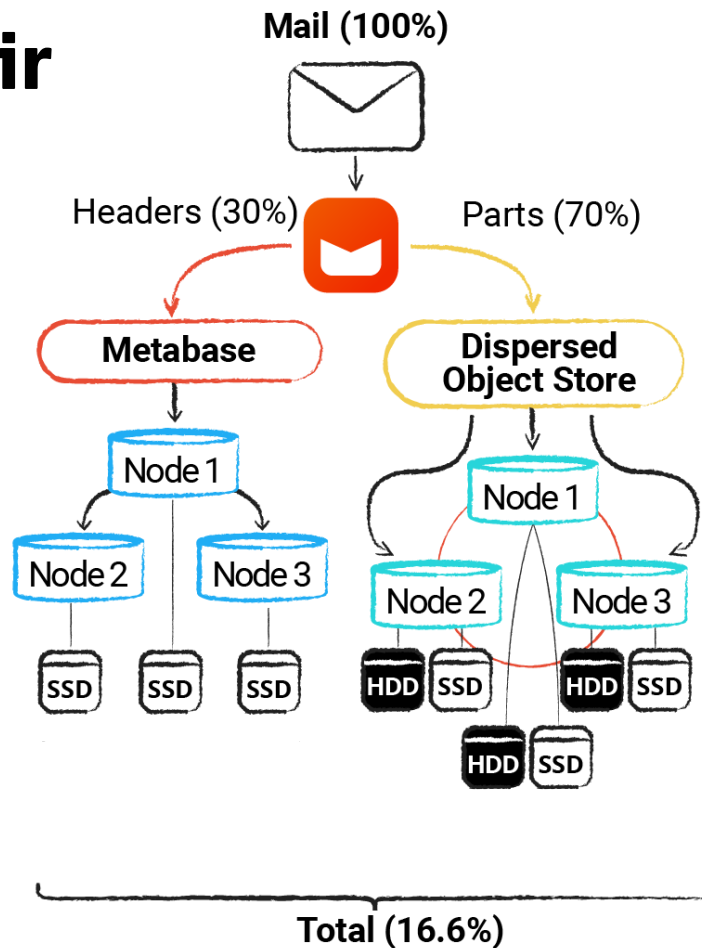


Mailion vs Maildir

$$RD_{meta} = 3$$

$$RD_{data} = \frac{2+2}{2} = 2$$

$$RD = 3$$

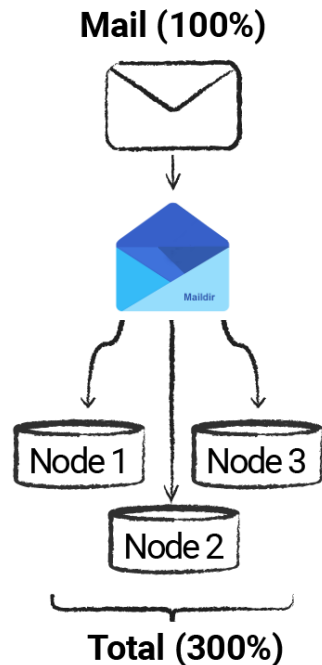
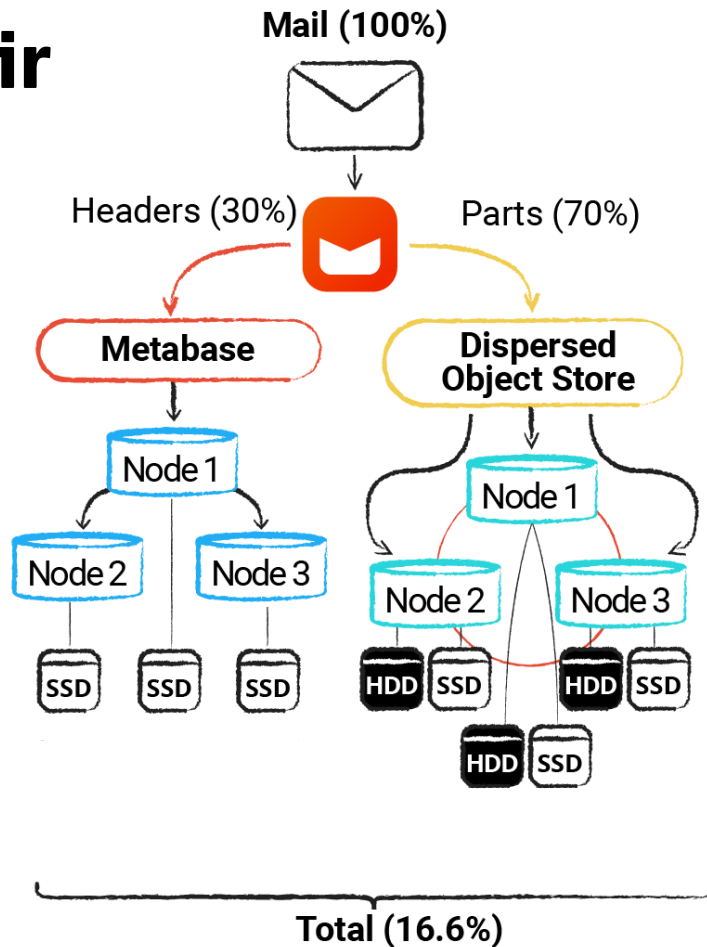


Mailion vs Maildir

- Потребление дискового пространства ниже в 18 раз (300% / 16.6% = 18.07).
- Стоимость хранения данных ниже в 8 раз.

$$RD_{meta} = 3$$
$$RD_{data} = \frac{2+2}{2} = 2$$

$$RD = 3$$



Расчёт стоимости хранения данных

Wikibon HDD & SSD Price-reduction Projections 2020 - 2030												
	2020	2021	2022	2023	2024	2025	2026	2027	2028	2029	2030	CAGR 2020-30
HDD \$/TB	\$ 22	\$ 21	\$ 19	\$ 18	\$ 17	\$ 16	\$ 15	\$ 14	\$ 14	\$ 13	\$ 13	-5.6%
SSD \$/TB	\$ 128	\$ 86	\$ 58	\$ 40	\$ 29	\$ 21	\$ 15	\$ 12	\$ 9	\$ 7	\$ 6	-26.8%
Ratio SSD/HDD	5.8	4.2	3.0	2.2	1.7	1.3	1.0	0.8	0.7	0.54	0.45	-22.5%
HDD Yr % decrease	-7.2%	-6.9%	-6.5%	-6.2%	-5.9%	-5.6%	-5.4%	-5.1%	-4.9%	-4.6%	-4.4%	
SSD Yr % decrease	-34%	-33%	-32%	-31%	-29%	-28%	-26%	-25%	-23%	-22%	-20%	
Ratio SSD/HDD	-29%	28%	27%	26%	25%	23%	22%	20%	19%	18%	16%	

Source: © Wikibon, 2021

- **500TB x 3 = 1,5PB (cluster Maildir)**
- **1,5PB x 19\$ = 28 500\$ (cluster Maildir)**

<https://wikibon.com/qlc-flash-hamrs-hdd>

Расчёт стоимости хранения данных

Wikibon HDD & SSD Price-reduction Projections 2020 - 2030												
	2020	2021	2022	2023	2024	2025	2026	2027	2028	2029	2030	CAGR 2020-30
HDD \$/TB	\$ 22	\$ 21	\$ 19	\$ 18	\$ 17	\$ 16	\$ 15	\$ 14	\$ 14	\$ 13	\$ 13	-5.6%
SSD \$/TB	\$ 128	\$ 86	\$ 58	\$ 40	\$ 29	\$ 21	\$ 15	\$ 12	\$ 9	\$ 7	\$ 6	-26.8%
Ratio SSD/HDD	5.8	4.2	3.0	2.2	1.7	1.3	1.0	0.8	0.7	0.54	0.45	-22.5%
HDD Yr % decrease	-7.2%	-6.9%	-6.5%	-6.2%	-5.9%	-5.6%	-5.4%	-5.1%	-4.9%	-4.6%	-4.4%	
SSD Yr % decrease	-34%	-33%	-32%	-31%	-29%	-28%	-26%	-25%	-23%	-22%	-20%	
Ratio SSD/HDD	-29%	28%	27%	26%	25%	23%	22%	20%	19%	18%	16%	

Source: © Wikibon, 2021

- **500TB x 3 = 1,5PB (cluster Maildir)**
- **1,5PB x 19\$ = 28 500\$ (cluster Maildir)**
- **28 500\$ / 8 = 3 562,5\$ (cluster DOS)**

<https://wikibon.com/qlc-flash-hamrs-hdd>

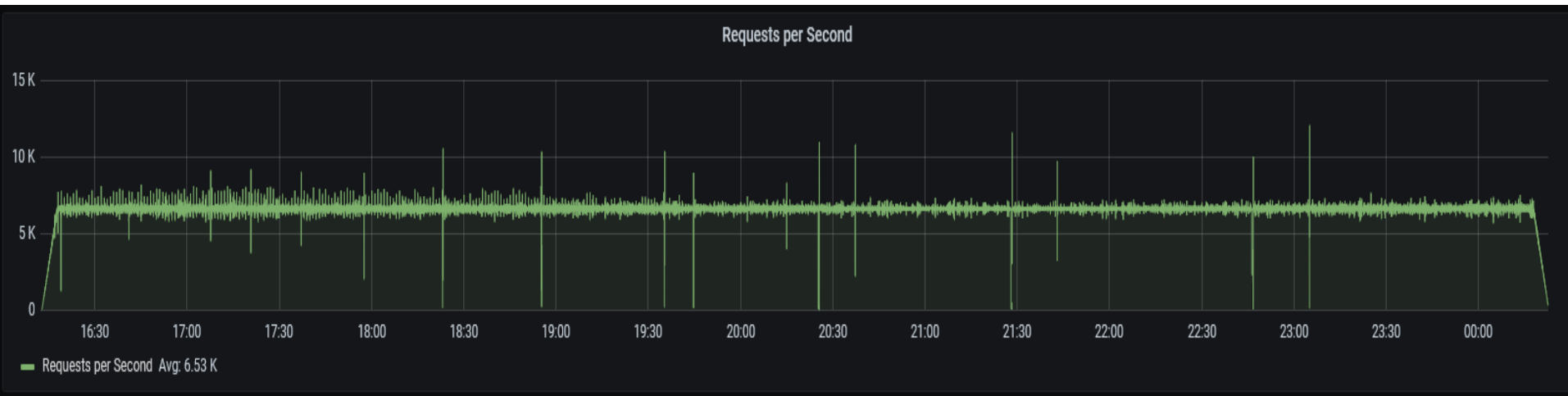
Нагрузочное тестирование

- Профиль нагрузки на 600 000 пользователей:
 - вход в систему;
 - отправка новых писем;
 - создание событий в календаре и реагирование на них;
 - другое.
- Для тестирования использовали K6 компании Grafana Labs
- Скрипты тестирования запускались на группировке из 46 серверов, которые суммарно были оснащены 636 процессорными ядрами, 2,8 ТБ оперативной памяти и накопителями емкостью более 135 ТБ.



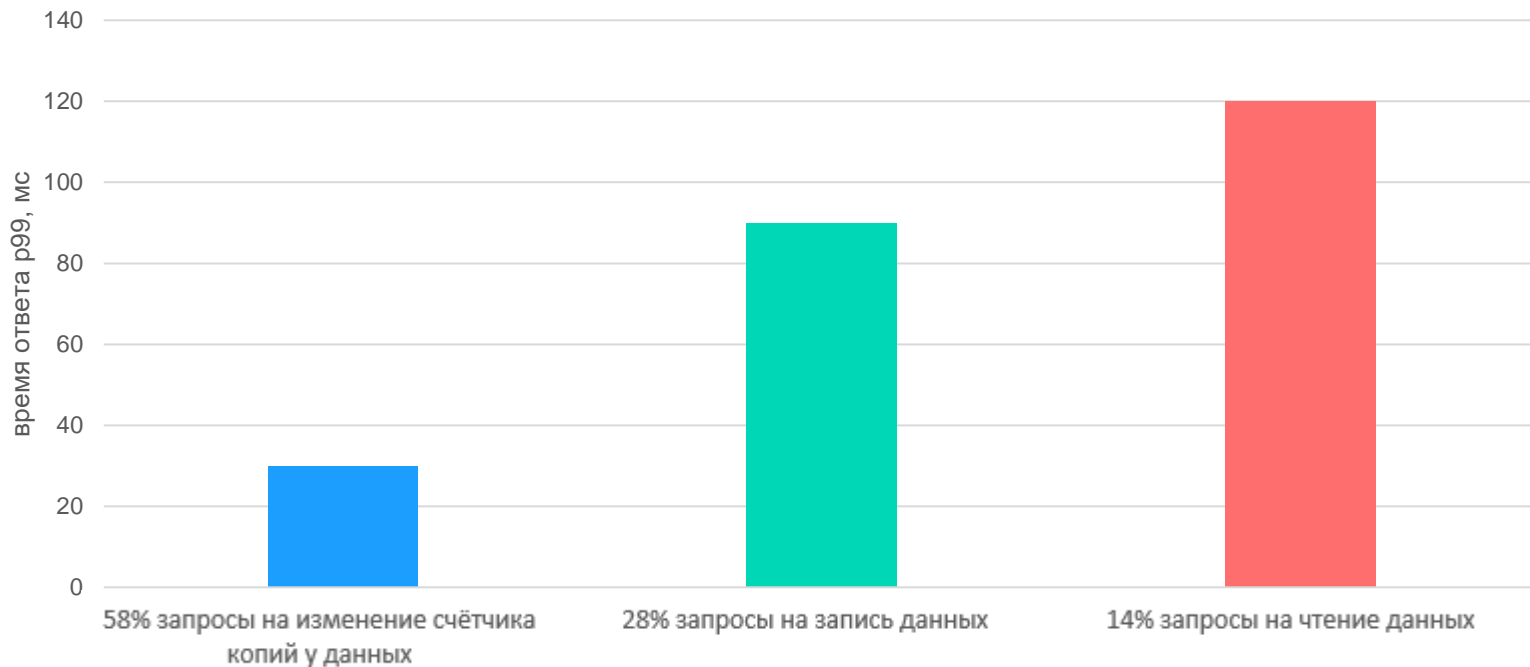
Нагрузочное тестирование

- В ходе испытаний инженерами проверялась гипотеза о стабильной работе системы под нагрузкой **6081 RPS** (операций в секунду). Это эквивалентно действиям **600 000 пользователей**, которые в течение дня отправляют и получают **1,14 млн писем**.



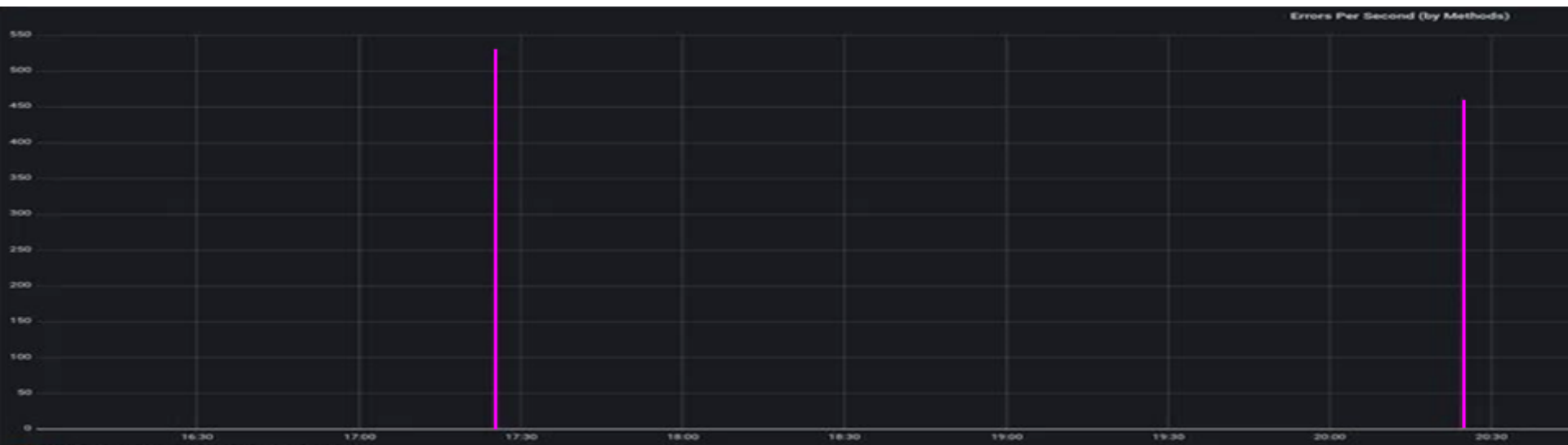
Нагрузочное тестирование: 6081 RPS

Распределение нагрузки на DOS



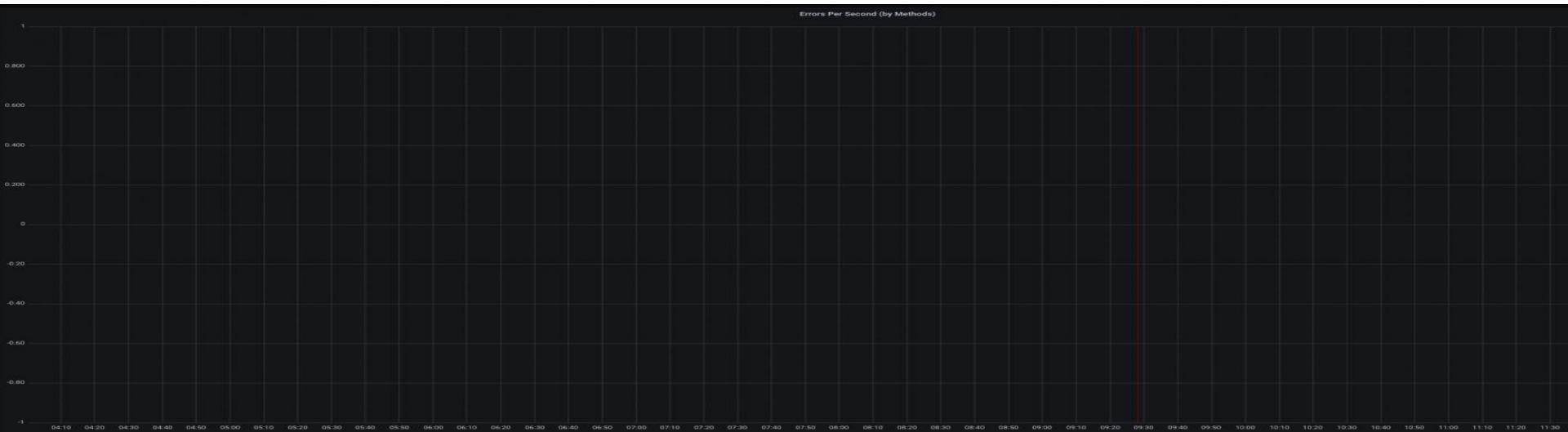
Нагрузочное тестирование: раунд 1

- Выявлены критические задержки некоторых методов



Нагрузочное тестирование: раунд 2

- Корректировка конфигурации (подключения к базе, кэш)



Собственное объектное хранилище



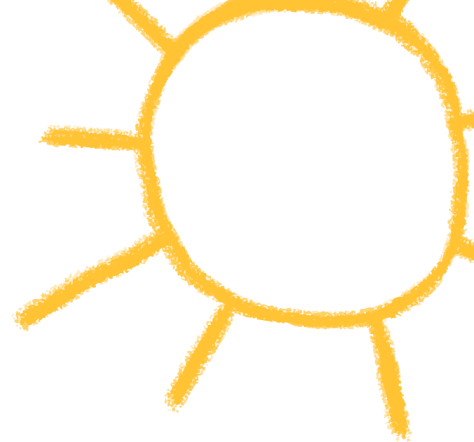
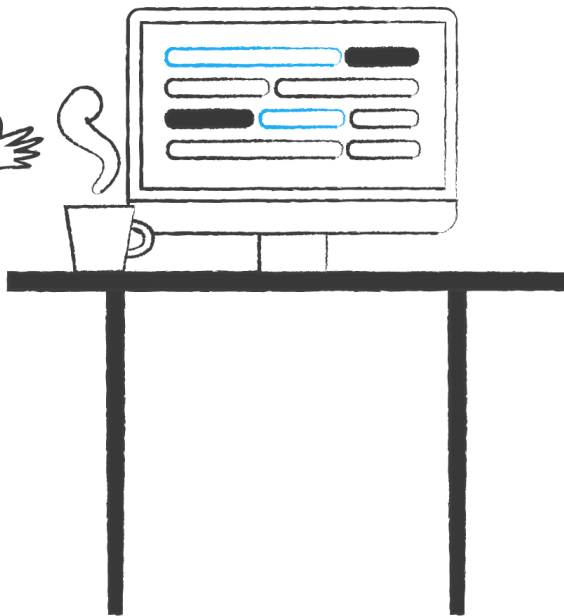
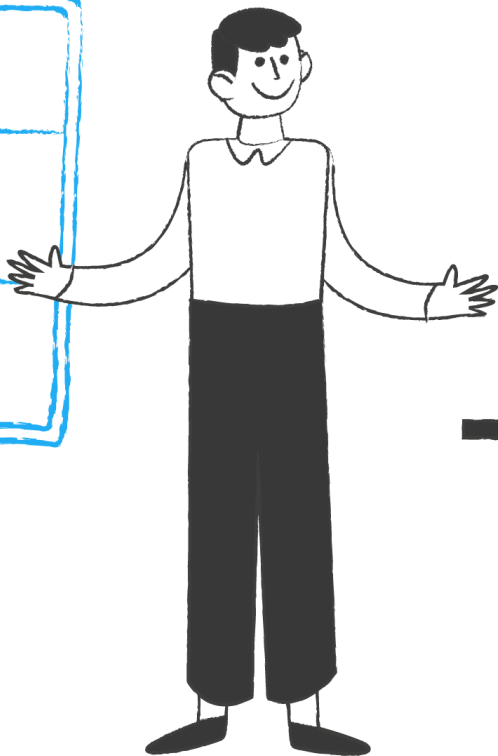
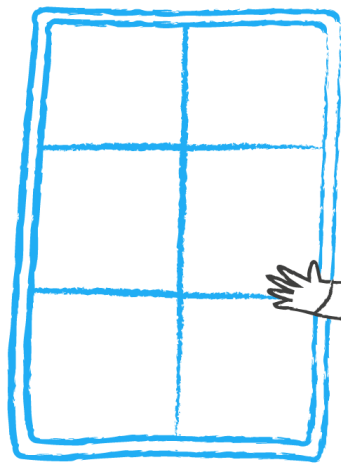
Плюсы:

- Контроль над разработкой, приоритизацией
- Дедупликация и компрессия
- Повышение экономической эффективности почтовой системы
- Оптимальный баланс аппаратных ресурсов (CPU / RAM vs Disk / IO)



Минусы:

- Сложная и дорогая разработка (распределённый stateful-сервис)



МойОфис



Андрей Колесников

☎ 8-929-384-19-24

✉ av_kolesnikov



Оценить доклад



HighLoad⁺⁺
2022

МойОфис

Слайды: <https://bit.ly/3sM0ihW>